



Determine the Most Effective Machine Learning Technique for Detecting Phishing Websites

Sm Mahamudul Hasan, Nirjas Mohammad Jakilim and
Md Forhad Rabbi

EasyChair preprints are intended for rapid
dissemination of research results and are
integrated with the rest of EasyChair.

August 31, 2021

Determine the Most Effective Machine Learning Technique for Detecting Phishing Websites

SM Mahamudul Hasan¹[0000-0002-4641-9451], Nirjas Mohammad Jakilim¹[0000-0003-3474-4908], and Md Forhad Rabbi¹

Shahjalal University of Science and Technology, Sylhet, Bangladesh
`smmahamudul32@student.sust.edu`

Abstract. Consumer tastes have moved away from conventional shopping and toward electronic commerce due to the Internet's fast growth. Rather than conducting bank or shop robberies, today's criminals use a range of sophisticated cyber methods to track down their victims. Attackers have developed new ways of deceiving customers, such as phishing, using fake websites to gather sensitive information such as account IDs, usernames, and passwords. The semantic-based nature of the assaults, which mainly leverage the vulnerabilities of computer users, makes establishing the authenticity of a web page more difficult. Machine learning (ML) is a typical data analysis technique that has shown promise in the battle against phishing. The article examines the applicability of machine learning methods for identifying phishing attempts and their advantages and disadvantages. Specifically, a variety of machine learning methods have been explored to find appropriate anti-Phishing technology solutions. More significantly, we put a wide range of machine learning methods to test real-world phishing datasets and against several criteria. To detect phishing websites, six different machine learning classification methods are employed. The Random Forest classifier had the most outstanding possible accuracy of 97.17% in this research, while the Gradient Boost Classifier had the highest achievable accuracy of 94.75%. The Decision Tree classifier has a provisioning accuracy of 94.69%. In contrast, Logistic Regression has a provisioning accuracy of 92.76%, KNN has a provisioning accuracy of 60.45%, and SVM has 56.04%. We showed that KNN has trouble detecting phishing sites since it hasn't been updated for accuracy. Decision trees are almost similar to Gradient Boosting in terms of performance.

Keywords: Phishing Detection, Website Security , machine learning classification, Random Forest, Decision Tree, SVM , Gradient Boost classifiers.

1 Introduction

1.1 Context

Phishing is one of the most serious cybersecurity threats since it includes creating fake websites that seem to be legitimate.[1] In this assault, the user inputs

important information such as credit card numbers, passwords, and so on to a bogus website that seems to be real. The sectors most severely affected by this attack include online payment services, e-commerce, and social media.[18][19] Phishing attacks make use of the aesthetic resemblance between fake and legitimate websites.[20] Phishing has taken on several forms throughout history, including legal, educational, and awareness efforts [21]. Phishing attacks use several methods to access sensitive information, including link manipulation, filter evasion, website forgery, covert redirection, and social engineering.[2,3] It is estimated that internet-based theft, fraud, and exploitation would account for an astonishing \$4.2 billion in financial losses in 2020, according to the FBI’s Internet Crime Complaint Center 2020 report.[4]

The attacker produces a website that is identical to the genuine web page in appearance. The Phishing web page’s URL is subsequently sent to thousands of Internet users through email and other contact forms[22]. Typically, the false email content creates panic, urgency, or promises money in exchange for the recipient taking immediate action. The fake email will prompt customers to change their PIN to prevent their debit/credit card suspension. When a user changes their sensitive credentials inadvertently, cyber thieves get the user’s information.[5] Phishing attacks are not just used to get information, they have also become the primary technique for distributing other kinds of harmful software, such as ransomware. 91% of current cyber-attacks begin with Phishing emails.[6]

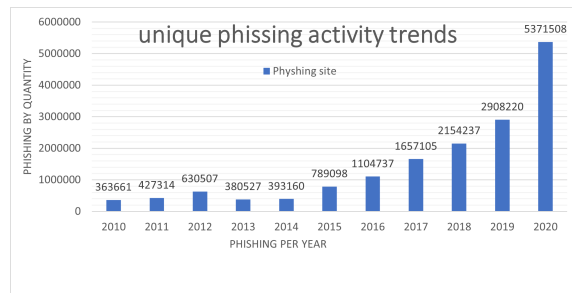


Fig. 1: Unique Phishing activity Trends

Phishing assaults account for more than half of all cyber fraud affecting Internet users. More than **245,771** distinct Phishing websites were identified in January 2021, according to the **APWG** study. Monthly attack growth rose by **1477%** during a ten-year period from 2010 to 2020. (363661 Phishing attacks in 2010 and an average of **5371508 attacks in 2020**). From 2010 through 2020, **Fig. 1** depicts the rise of Phishing assaults.[7]

1.2 Problem definition

Phishing is the act of creating and duplicating legitimate websites with the intent of duping internet users into disclosing their login credentials and personal data[23]. Numerous variations have occurred throughout history. The attacker may be motivated by identity theft, financial gain, or celebrity. Scientists and researchers have significant challenges when it comes to identifying and blocking Phishing attacks[24]. Even the most seasoned and informed users may face an assault. Phishing attacks usually include the transmission of a fake email purporting to be from a reputable company or organization, requesting sensitive information such as a bank login or password. Phishing communications are delivered through email, SMS, instant messaging, social media, and voice over internet protocol (VoIP).[8] However, the most frequent form of attack is through email. 65% of Phishing efforts include the use of malicious URLs in emails.[9]

Due to the fact that Phishing websites are only up for two days, the phisher may flee immediately after committing the crime. While user education may aid in the worldwide battle against Phishing, it is expensive and needs users to be familiar with computer security.

1.3 Proposed Solution

For determining if a website is phishing, a variety of methods have been developed. Attackers may use several methods at various phases of the attack cycle. Among other things, network security, user education, user authentication, server-side filters, client-side tools, and classifiers are some of the methods that are available. While each kind of Phishing attack is unique, the bulk of them have certain characteristics and patterns.[10] Due to the essential role of machine learning techniques for identifying patterns in data, it has become possible to identify a large number of common Phishing characteristics, as well as to recognize Phishing websites.[11] Specifically, the purpose of this paper is to examine and evaluate a number of machine learning methods for identifying fake websites. Logistic Regression, Random Forest, Support Vector Machine, KNN, Naive Bayes, and the XGBoost classifier were among the machine learning methods we investigated.

1.4 Experimental result

We have tested our suggested Phishing detection method on different classification algorithms and utilized a dataset of **11054** Phishing and non-Phishing websites. Experimental findings indicate that logistic regression works best in the identification of Phishing websites. The suggested method has reasonably high accuracy in identifying Phishing websites as it was obtained for Random forest classification. It gives more than 97.41% true positive rate and just 2.58% false positive rate. Moreover, our method's accuracy, precision, and f1 score are 97.17%, 97.80, and 97.59%, respectively. We have also examined the area under

the classification model learning curve for all the models to discover a better measure of accuracy. Our experiment calculated training and cross-validation scores independently for all classification models used to categorize correct web pages.

1.5 Outlines

The remainder of this research paper is divided into the following sections: The second section covers similar studies on phishing websites. Section 3 explored into our suggested method of work. In Section 4, we provide a short description of the dataset. In Section 5, we examine different machine learning techniques for detecting Phishing and summarized the Experiments and its outcome. The conclusion and future study are described in Section 6.

2 RELATED WORK

There are methods for phishing detection that are list-based and machine-learning-based. The most often used detection technique is list-based. Whitelists consist of legitimate websites, while blacklists consist of phishing websites. Phishing detection systems that are list-based rely on these lists to identify phishing attempts. Whitelists enable the development of secure, authentic, and comprehensive websites. Unwhitelisted websites may be hazardous. C Whittaker, B Ryner, M Nazif[25] compiled and published a whitelist of all URLs accessed through the Login user interface for. When a user visits a website, the system informs the user if their data is incompatible with the site. This method may be used by a user the first time they visit an authorised website.

SL Pflieger, G Bloom[26] developed an automatic whitelist of user-approved websites. Two stages of feature extraction are domain-IP address matching and source code connections. True positives accounted for 86.02% of all observations, whereas false negatives accounted for 1.48%.

Zhang et al. [12] identified phishing on behalf of CANTINA by using the frequency-inverse document frequency technique. After that, the keywords were put into Google. When a website appears in search results, it earns the confidence of users. CANTINA Plus is equipped with fifteen HTML-based functionalities (Researchers4). Despite the algorithm's 92% accuracy rating, a substantial number of false positives may have happened.

The title and message were reviewed by Islam et al.[13] A categorization system for communication was created. The technique is designed to minimise false positives. The research gathered information on URL-specific characteristics such as length, subdomain names, slashes, and dots. Rule mining was utilised to develop detection rules as a priority. In testing, 93% of phishing URLs were identified.

To categorise phishing attempts more correctly and effectively, an adaptive self-structuring neural network was employed by Rami et al.[14] It includes seventeen features, some of which are dependent on third-party services. As a result,

real-time execution is sluggish. While it has the potential to enhance accuracy, it is not presently used. It manages noisy data by using a small dataset of 1400 elements.

Others combine artificial intelligence and image processing. For image/visual-based applications, the internet domain (web history) is needed. Recognize phishing attempts (web history).[15] The proposed approach circumvents these constraints. They categorised features according to whether they included hyperlinks, third-party material, or masked URLs. By using third-party services, the system's accuracy is increased to 99.55%, while detection time and latency are decreased.

NLP receives little attention in the scholarly literature (NLP). A recent research Peng et al.[16] used natural language processing to detect phishing emails. This programme scans the plain text content of emails for harmful intent. It gathers queries and responses via the use of natural language processing (NLP). Phishing attempts are detected using a custom-built blacklist of word pairs. The algorithm was trained on 5009 phishing emails and 5000 real emails prior to becoming public. Their experimental research demonstrated a 95% accuracy rate.

3 PROPOSED APPROACH

The method we use to identifying Phishing websites is one that is based on machine learning. Our model incorporates a variety of techniques, including logistic regression, KNN, decision trees, Random Forest, support vector machines, and gradient boosting. We were supplied with the whole dataset via Kaggle.[17]

The dataset includes **11054** records, including information about Phishing and genuine websites, as well as **32 features**. This dataset was generated via the extraction of source code and the construction of a DOM tree. Then, with the use of a web browser, we extract 32 characteristics from specific websites. These characteristics determine whether or not a website is Phishing. We gather data from Kaggle [17] and highlight the dataset's vector in our model. Then, we utilized this dataset for training six machine learning classification models to identify the characteristics of a Phishing website.

Our categorization algorithm detects Phishing websites with an accuracy of about 97%. And our model is capable of detecting approximately 97.48% of genuine positive confusion matrixes. We successfully anticipate Phishing websites using our machine learning classification algorithm. This is our general strategy for detecting Phishing websites.

4 DATASET

One of the most important challenges we faced during our research was a scarcity of Phishing databases. There have been a large number of academic papers on the subject of Phishing detection; however, none of them have made the datasets used in their research available to the general public. The absence of a common

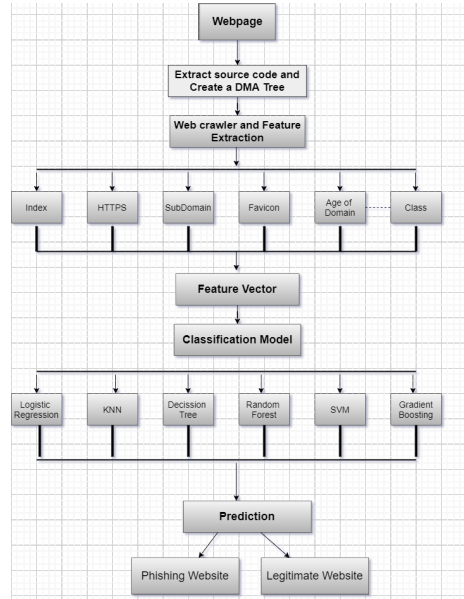


Fig. 2: Diagram of Our Approach

feature set that captures the features of a Phishing website also makes it more difficult to collect useful data, which makes it more difficult to build a useable dataset in the first place. Numerous academics carefully examined and benchmarked the dataset used in our analysis. We collect our dataset from kaggle.[17] The collection contains about **11,054** sample websites with 32 features. Nearly **6080** legitimate websites and **4974** Phishing websites are included in our collection. In our dataset, we give a score of 1 to genuine websites and -1 to Phishing websites. We utilized 30% of samples for testing and 70% for training. Each website is evaluated to see if it is legitimate or fake.

We categorize the **32 features** in our dataset into three categories. The following categories are described:

4.1 Address bar based features

This element is considered in the address bar-based features. Long URL, Short URL, Symbol@, Redirecting/, Prefix Suffix-, Subdomains, HTTPS, DomainRegLen, Favicon, NonStdPort, HTTPSDomainURL, Request URL, and Anchor URL are all used in the index. These **15 features** are together referred to as address bar-based functionalities. The address bar is often referred to as the link URL of features. The address bar is described in the following manner:



Fig. 3: URL the component of a legitimate website.

4.2 Abnormal based features

There are **9 features** that define abnormal based characteristics. LinksInScript-Tags, Server Form Handler, Info Email, Abnormal URL, Website Forwarding, StatusBarCust, Disable Right Click, Using popup Window, and IframeRedirection are some of them.

4.3 Domain based features

There are **8 features** that describe characteristics that are abnormal. These are the following: Domain Age, DNSRecording, Website Traffic, Links Pointing to Page, PageRank, Google Index, Status Report, and class.

If the URL feature has a value of -1, it is a phishing website. If the URL value is 0, the website is suspect. If the URL value includes one, it indicates that the website is genuine.

5 RESULT AND ANALYSIS

We utilized 10-fold cross-validation to evaluate the model’s overall performance in our trials. 10 sub-samples were drawn from the original data set using a random number generator. Three samples are tested (30%), while the other samples are used to train model-based categorization algorithms. Because phishing detection is categorical in nature, we must employ a binary classification model to identify phishing assaults. ”-1” denotes a phishing sample, while ”1” denotes a genuine sample. We identified phishing websites using a variety of machine learning models, including logistic regression, random forest, KNN, SVM, gradient boosting, and decision trees.

We assessed these models’ accuracy, precision, recall, F1 score, and confusion matrix, and then utilized a variety of feature selection and hyperparameter tweaking techniques to get the best possible results. The precision, recall, and F1 scores and accuracy of different models, as well as their overall performance, are summarized in **table 1**. The accuracy of different models, as well as their overall performance, is compared in **fig 5**. The Random Forest has been proven

to be extremely accurate, reasonably resistant to noise and outliers, simple to build and comprehend, and capable of implicit feature selection in our studies. In the fig 4, we show the learning curves of random forest. The Random Forest offers a number of benefits over the decision tree, the most important being its resistance to noise.

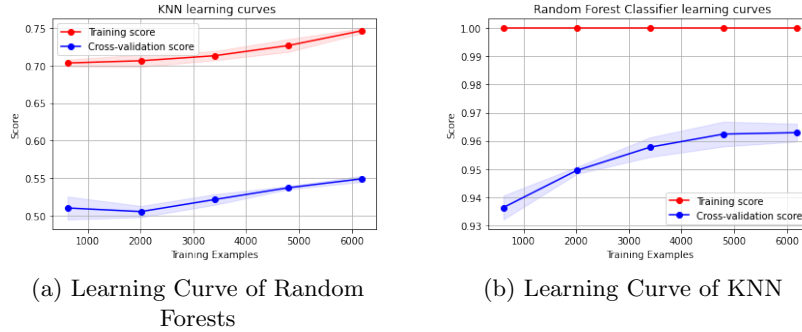


Fig. 4: Learning Curve of RF and KNN

By increasing the number of trees in each woodland inside the forest, the Random Woodland decreases variance. The primary drawback of Random Forests was the large number of hyperparameters that required tuning to attain optimum performance. Additionally, it adds a random aspect to both the training and testing data, which may not be appropriate for all data sets and circumstances. The study could establish the optimal classification accuracy of the KNN by using $k=10$. There is no one-size-fits-all k value for KNN classification. In the fig 4, we show the learning curves of KNN classifier. Due to the huge number of neighbours, it is computationally costly to develop a solution. Additionally, we discovered that a few neighbours provide the most flexible fit, with low bias but high variation, while a large number of neighbours produces a smoother decision boundary with lower variance but greater bias.

Logistic regression is predicted to be 92.76% accurate. Additionally, our system accurately detects about 93.50% of true positives in confusion matrix. We can evaluate the model's correctness throughout both the training and testing stages by examining the training and cross validation scores. The actual positive accuracy of the KNN model is just 55.28%, and the model can not identify 44.71% of phishing websites. This accuracy is poor, and the total performance of the KNN model is 60.45%. When the accuracy of decision tree tests is considered, the cross-validation score performs well. The decision tree classifier has a true positive score of 95.32%, but a false positive score of just 4.67%. With an overall accuracy of 94.69%, the model is very accurate. The Random forest has the greatest accuracy among all the models. At level 1, the training score of the model is negligible. Cross-validation surpasses single-validation. This model has

a true positive rate of 97.48% and a false negative rate of 2.51%. The system works optimally 97.17% of the time. The support vector machine model performs the least well of all the models. Accuracy is optimal 56.0% of the time. It is unable to identify any phishing websites effectively. The training score is the lowest, but the cross-validation score is close to the training score overall. As a result, the model is unable to function properly. In our model, the gradient boosting method works well. It detects true positives at a rate of 95.53% and false positives at a rate of just 4.46%. The model’s total accuracy is 94.75%.

Table 1: Evaluation of all the models

Algorithms .	Precision .	Recall .	F1score .	Accuracy
Logistic Regression	0.92	0.93	0.93	92.76%
KNN	0.52	0.55	0.54	60.45%
Decission Tree	0.95	0.95	0.95	94.75%
Random Forest	0.97	0.96	0.97	97.17%
SVM	0.50	0.28	0.36	56.04%
Gradient Boosting	0.94	0.95	0.95	94.75%

The primary benefit of Gradient Boost over other techniques such as decision trees and support vector machines are its speed. Additionally, it has a regularization parameter that significantly lowers variance. To further enhance the generalizability of this approach, the learning rate, and subsamples from features such as random forests are combined with the regularization parameter. Compared to Logistic Regression and Random Forests, Gradient Boost is more difficult to comprehend, visualize, and change. Numerous hyperparameters may be adjusted to improve overall performance. Gradient Boost is an enticing technique to use when both speed and accuracy are required. Despite this, more resources are needed to train the model, since model tweaking takes additional effort and skill on the part of the user to get statistically significant results.

The decision tree outperforms the KNN, Logistic regression, and SVM in our model. Due to the huge volume of data and the diversity of characteristics included therein, the decision tree works well in this scenario. The Decision Tree has two nodes. The Decision Node is the first of these nodes, followed by the Leaf Node. In contrast to decision nodes, which are used to make choices and include many branches, leaf nodes reflect the result of those choices and contain no further branches that branch out to other locations. The judgments or tests are done in light of the dataset’s characteristics.

For SVM, we just apply the linear kernel model. Our prior experience indicates that the linear kernel does not perform well on this dataset. Consequently, SVM is ineffective at detecting phishing websites. It is unable to properly detect any phishing websites. Despite the size of our dataset, the SVM technique is not designed for big data sets. When the data set has a high level of noise and the target classes overlap, SVM performs poorly. It is unusual for the SVM to per-

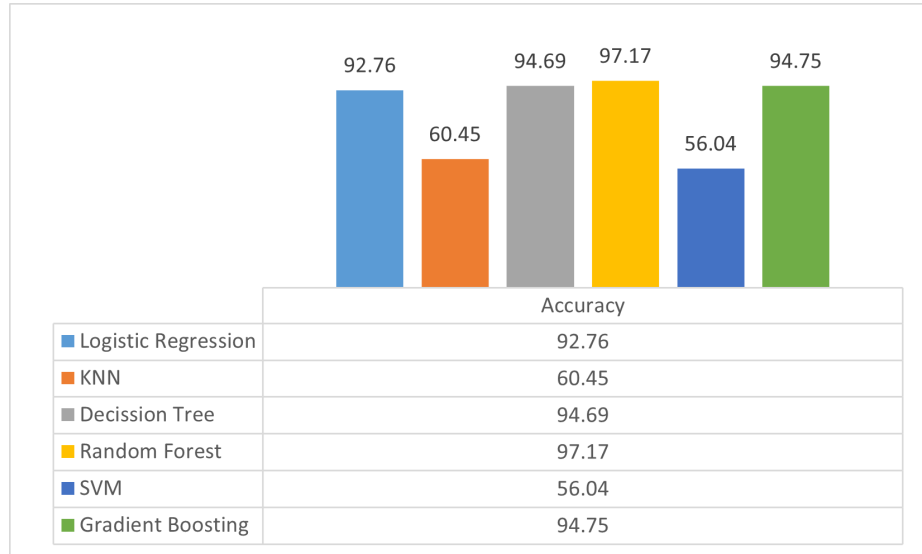


Fig. 5: Accuracies of the Models

form poorly when a single data point has more features than there are training data samples. Therefore, SVMs fail in our model.

6 CONCLUSION

We developed and tested six phishing website classifiers on a dataset of 6080 legitimate websites and 4974 phishing websites in this study. Classifiers such as Logistic Regression, Decision Trees, Support Vector Machines, Random Forests, KNNs, and Gradient Boosting are examined. Our classifiers, Random Forest and Gradient Boost, perform well in terms of computation time and accuracy, as shown in Tables 1. Experimental findings indicate that logistic regression works best for the identification of phishing websites. The suggested method has reasonably high accuracy in identifying phishing websites as it was obtained for random forest classification. It is more than 97.41% true positive rate and just a 2.58% false positive rate. Moreover, our method's accuracy, precision, and f1 score are 97.17%, 97.80, and 97.59%, respectively. We have also examined the area under the classification model, learning curve for all the models to discover a better measure of accuracy. Our experiment calculated training and cross-validation scores independently for all classification models used to categorize correct web pages.

Our model could not make use of support vector machines because SVM uses the linear kernel model. That's why When the input is noisy and the target classes overlap, SVM performs poorly. When a single data point has more features than the amount of training data samples, it fails. As a consequence, SVMs

don't perform optimally in our model. Utilizing a different SVM kernel may be beneficial. Also, using a polynomial, Sigmoid, or RBF kernel may increase accuracy. Logistic regression is predicted to be 92.76% accurate. The KNN accuracy is poor, and the total performance of the KNN model is 60.45%. When the accuracy of decision tree tests is considered, the cross-validation score performs well. The decision tree classifier has a true positive score of 95.32%, but a false positive score of just 4.67%. With an overall accuracy of 94.69%, the model is very accurate.

References

1. Shaikh, Anjum N., Antesar M. Shabut, and M. Alamgir Hossain. "A literature review on phishing crime, prevention review and investigation of gaps." 2016 10th International Conference on Software, Knowledge, Information Management & Applications (SKIMA). IEEE, 2016.
2. Ali, Abdul. "Social Engineering: Phishing latest and future techniques." Retrieved March 10 (2015): 2016.
3. Goel, Diksha, and Ankit Kumar Jain. "Mobile phishing attacks and defence mechanisms: State of art and open research challenges." *Computers & Security* 73 (2018): 519-544.
4. FBI Releases the Internet Crime Complaint Center 2020 Internet Crime Report, Including COVID-19 Scam Statistics, <https://www.fbi.gov/news/pressrel/press-releases/fbi-releases-the-interne-crime-complaint-center-2020-internet-crime-report-including-covid-19-scam-statistics>
5. Dhamija, Rachna, J. Doug Tygar, and Marti Hearst. "Why phishing works." Proceedings of the SIGCHI conference on Human Factors in computing systems. 2006.
6. 91% of all cyber attacks begin with a phishing email to an unexpected victim, <https://www2.deloitte.com/my/en/pages/risk/articles/91-percent-of-all-cyber-attacks-begin-with-a-phishing-email-to-an-unexpected-victim.html>
7. PHISHING ACTIVITY TRENDS REPORTS, <https://apwg.org/trendsreports/>
8. Ramzan, Zulfikar. "Phishing attacks and countermeasures." Handbook of information and communication security (2010): 433-448.
9. Must-Know Phishing Statistics: Updated 2021, <https://www.tessian.com/blog/phishing-statistics-2020/>
10. Jain, Ankit Kumar, and B. B. Gupta. "A survey of phishing attack techniques, defence mechanisms and open research challenges." *Enterprise Information Systems* (2021): 1-39.
11. Passos, Ives Cavalcante, Benson Mwangi, and Flávio Kapczinski. "Big data analytics and machine learning: 2015 and beyond." *The Lancet Psychiatry* 3.1 (2016): 13-15.
12. Zhang, Yue, Jason I. Hong, and Lorrie F. Cranor. "Cantina: a content-based approach to detecting phishing web sites." Proceedings of the 16th international conference on World Wide Web. 2007.
13. Islam, Rafiqul, and Jemal Abawajy. "A multi-tier phishing detection and filtering approach." *Journal of Network and Computer Applications* 36.1 (2013): 324-335.

14. Mohammad, Rami M., Fadi Thabtah, and Lee McCluskey. "Predicting phishing websites based on self-structuring neural network." *Neural Computing and Applications* 25.2 (2014): 443-458.
15. Basit, Abdul, et al. "A comprehensive survey of AI-enabled phishing attacks detection techniques." *Telecommunication Systems* (2020): 1-16.
16. Peng, Tianrui, Ian Harris, and Yuki Sawa. "Detecting phishing attacks using natural language processing and machine learning." 2018 IEEE 12th International Conference on Semantic Computing (ICSC). IEEE, 2018.
17. Phishing Website Detector,
<https://www.kaggle.com/eswarchandt/phishing-website-detector>
18. Scheau C., A. Arsene, and Gerald Dinca. "Phishing and e-commerce: an information security management problem." *Journal of Defence Resources Management* 7.1 (2016): 12.
19. Sarjiyus, O., N. D. Oye, and B. Y. Baha. "Improved Online Security Framework for e-Banking Services in Nigeria: A Real World Perspective." *Journal of Scientific Research and Reports* (2019): 1-14.
20. Mohammad, Rami M., Fadi Thabtah, and Lee McCluskey. "Tutorial and critical analysis of phishing websites methods." *Computer Science Review* 17 (2015): 1-24.
21. Adebowale, Moruf A., et al. "Intelligent web-phishing detection and protection scheme using integrated features of Images, frames and text." *Expert Systems with Applications* 115 (2019): 300-313.
22. Jain, Ankit Kumar, and Brij B. Gupta. "A novel approach to protect against phishing attacks at client side using auto-updated white-list." *EURASIP Journal on Information Security* 2016.1 (2016): 1-11.
23. Charoen, Danuvasin. "Phishing: a Field Experiment." *International Journal of Computer Science and Security (IJCSS)* 5.2 (2011): 277.
24. Jakobsson, Markus, and Steven Myers, eds. *Phishing and countermeasures. understanding the increasing problem of electronic identity theft*. John Wiley & Sons, 2006.
25. @articlewhittaker2010large, title=Large-scale automatic classification of phishing pages, author=Whittaker, Colin and Ryner, Brian and Nazif, Marria, year=2010
26. Pfleeger, Shari Lawrence, and Gabrielle Bloom. "Canning spam: Proposed solutions to unwanted email." *IEEE Security & Privacy* 3.2 (2005): 40-47.